



(12) **EUROPEAN PATENT APPLICATION**

(43) Date of publication:
30.10.2002 Bulletin 2002/44

(51) Int Cl.⁷: **G06F 9/50**

(21) Application number: **02009207.8**

(22) Date of filing: **24.04.2002**

(84) Designated Contracting States:
AT BE CH CY DE DK ES FI FR GB GR IE IT LI LU
MC NL PT SE TR
 Designated Extension States:
AL LT LV MK RO SI

(72) Inventors:
 • **Dorofeev, Andrei V.**
Sunnyvale, CA 94086 (US)
 • **Tucker, Andrew G.**
Menlo Park, CA 94025 (US)

(30) Priority: **25.04.2001 US 843426**

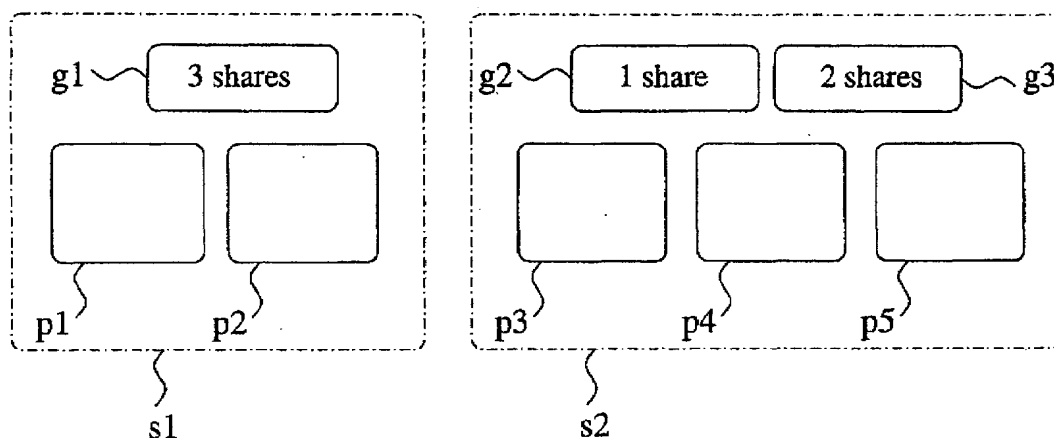
(74) Representative: **HOFFMANN - EITLE**
Patent- und Rechtsanwälte
Arabellastrasse 4
81925 München (DE)

(71) Applicant: **Sun Microsystems, Inc.**
Palo Alto, California 94303 (US)

(54) **Apparatus and method for scheduling processes on a fair share basis**

(57) Described is a scheduling system that provides allocation of system resources of one or more processor sets among groups of processes. Each of the process groups is assigned a fixed number of shares, which is the number that is used to allocate system resources among processes of various process groups within a given processor set. The described fair share scheduler considers each processor set to be a separate virtual computer. Different process sets do not share processes, a particular process must execute on a single processor set.

In another embodiment of the invention, each process group could be given a separate number of shares for each processor set. Percentage of the resources of the specific processor set allocated to processes of a process group is calculated as a ratio of the shares of the process group on the processor set to the total number of shares of active process groups operating in that set. The process group is considered active on a processor set, if that processor set executes at least one process of that process group.



$$P_{1-5} = 20\% \text{ (1/5th) of the system resources.}$$

Fig. 1

Description

FIELD

[0001] The present invention generally relates to techniques for managing load of a central processing unit in a computer system. The present invention also relates to techniques for allocating central processing unit percentage usage among several processes.

BACKGROUND

[0002] Modern high performance computer systems are capable of supporting multiple users and executing multiple processes or applications started by different users at the same time. Term scheduling refers to allocation of central processor unit (CPU) resources among multiple processes owned by different users in a computer system. It will be appreciated by persons of skill in the art that the fraction of the total CPU resources allocated to a particular process may depend on the number of other processes in the system and a relative importance of the process in comparison with the other processes in the system.

[0003] Another way of allocating system resources is to provide a process group, rather than a single process, with a fixed share of CPU resources. In other words, one or more processes executing on a central processing unit of a computer system can be combined into process groups, and each process group is allocated a share of the system resources that are distributed among individual processes of this process group. The system resources can be distributed among processes of the process group in an equal or unequal manner. The processes can be combined into aforementioned process groups based on various criteria, including, but not limited to, user id of the user executing the process, group id of the user executing the process, or based on any other appropriate classification. The nature of the classification is not critical to the present invention.

[0004] Modern operating systems include concept of processor sets, which is a group of processes executing a specific separate set of processes. Only those processes that are explicitly bound to a processor set are permitted to be executed on processors that belong to that processor set. Presently, there exists no resource allocation scheme that would leverage advantages of the processor set architecture and have an ability to distribute CPU resources among multiple process groups in a pre-defined manner.

SUMMARY

[0005] It is an object of the invention to provide for improved allocation of system resources.

[0006] According to an embodiment a method for allocating a percentage of system resources among process groups in a computer system, said computer system

comprising at least one central processing unit, said at least one central processing unit combined into at least one processor set, may comprise: a. assigning each of said process groups a number of shares for each or said at least one processor set; b. allocating said system resources of each of said at least one processor set to each of said process groups according to the number of shares assigned to said each of said process groups.

[0007] Said system resources of each of said at least one processor set may be allocated based on a number of shares of all active groups within said each of said at least one processor set.

[0008] Further, said percentage of said system resources may be calculated based on a ratio of the number of shares assigned to said each of said process groups to the a number of shares of all active groups within said each of said at least one processor set.

[0009] Still further, each of said process groups may include only one process.

[0010] A program may be provided having instructions adapted to make a computer carry out at least one of the above operations. And, a computer readable medium may be provided, embodying the program.

[0011] According to another embodiment a computer readable medium may be provided, embodying a program for allocating a percentage of system resources among process groups in a computer system, said computer system comprising at least one central processing unit, said at least one central processing unit combined into at least one processor set, said program comprising: a. assigning each of said process groups a number of shares for each or said at least one processor set; b. allocating said system resources of each of said at least one processor set to each of said process groups according to the number of shares assigned to said each of said process groups.

[0012] According to still another embodiment, a scheduler may be provide for allocating a percentage of system resources among process groups in a computer system with at least one central processing unit, said at least one central processing unit combined into at least one processor set, said scheduler comprising means for assigning each of said process groups a number of shares for each or said at least one processor set; and means for allocating said system resources of each of said at least one processor set to each of said process groups according to the number of shares assigned to said each of said process groups.

[0013] Further advantageous features of the invention are outlined in further claims.

DESCRIPTION OF THE DRAWINGS

[0014] Various embodiments of the present invention will now be described in detail by way of example only, and not by way of limitation, with reference to the attached drawings wherein identical or similar elements are designated with like numerals.

[0015] FIG. 1 illustrates exemplary allocation of system resources according to the inventive concept.

DETAILED DESCRIPTION

[0016] To overcome the limitations described above, and to overcome other limitations that will become apparent upon reading and understanding the present specification, apparatus, methods and articles of manufacture are disclosed that allocate a percentage of system resources among process groups in a computer system having one or more processor sets.

[0017] According to the inventive method, processes are combined into process groups based on a pre-defined criteria. Each of the process groups is assigned a number of shares representing relative importance of the group within its processor set. The inventive system allocates the system resources of the processor sets to process groups according to the number of shares assigned to a particular process group and the total number of shares of all active process groups in the processor set. A process group is considered active on a processor set if there is at least one process of this process group executing on that processor set.

[0018] Fair share scheduling is a way to assign a particular process a fixed share of CPU resources. The term share may be used to describe the relative importance of one workload versus another.

[0019] According to one of the aspects of the present invention, various processes in the system are combined into one or more process groups. These process groups users are assigned a number of shares which represent relative importance thereof. This is a way to guarantee application performance by explicitly allocating shares (or percentage) of system CPU resources among competing workloads. Note that total number of shares assigned to all process groups need not be 100. Furthermore, to obtain the percentage of the system CPU resources available to a process group at each given moment of time, a total number of shares allocated to that process group must be divided by the total number of shares possessed by all currently active process groups. A process group is considered active when it has at least one running or runnable process. Indeed, to ensure the complete, or 100% utilization of the system, only process groups which have executing processes at a particular time should be given share of the CPU usage. Note that such active process groups are searched across the entire system. At any given time, the percentage of the CPU allocated to a particular process group depends on the number of shares owned by all other active process groups in the system, or the process groups that have at least one executing process in the system. Therefore, in a system where processes are combined into process groups based on user id, any new logged user with a given number of shares can decrease the CPU percentage of all other actively running users.

[0020] Modern operating systems, such as a Solaris Operating System distributed by Sun Microsystems, Inc. of Palo Alto, California have a concept of processor sets. Processor set concept applies to multiple processor computers and allows the binding of one or more processors into groups of processors. Processors assigned to processor sets will run only processes that have been bound to that processor set. In other words, the aforementioned processor set is essentially a virtual single- or multi-processor computer system within a physical computer, which has its own set of running processes. The concept of processor set is especially helpful, for example, when certain important process need to be provided with a separate one or more processors. For example, in a computer system providing services to http clients, a separate processor set can be allocated to running a web server, while all other processes can be executed on a second, separate processor set. In this case, the amount of CPU resources allocated to the web server will not depend on the other processes executing in the system.

[0021] However, when the aforementioned processor sets are used in conjunction with the conventional fair share scheduler, the performance of processes running on one processor set may be impacted by the work performed by processes running on another processor set, which is an undesirable effect.

[0022] The reason why the existing fair share scheduler does not work satisfactorily with processor sets is because that total number of shares for all active process groups is calculates across the entire system, when in fact it should be calculated only within boundaries of the current processor set. If the total number of shares is kept separate for each processor set, then the CPU allocation for a given process group will only depend on other process groups who have their active processes on the same processor set. The work done on other processor sets will be unaffected. This is more intuitive behavior of such configurations than what it has been in the past.

[0023] The inventive fair share scheduler will now be described in detail. According to the inventive concept, various processes in a computer system are combined into process groups. Each of these process groups is assigned a fixed number of shares, which is the number that represents relative importance of process groups. The number of shares of a process group is used to allocate system resources among processes of that process group executing within a predetermined processor set, in the manner described in detail below. Specifically, the inventive fair share scheduler considers each processor set to be a separate virtual computer. Different processor sets do not share processes, in other words, a process must execute on a single processor set.

[0024] In one embodiment of the invention, each process group is given the same number of shares for all processor set. It should be noted that even if process group has zero shares, processes of this process group

may still be executed on processor sets, which do not have any other active process groups. Percentage of the resources of the specific processor set allocated to processes of a particular process group is calculated as a ratio of the shares of that process group on the processor set to the total number of shares of all active process groups operating in that set. The process group is considered active on a processor set, if that processor set executes at least one process of that process group.

[0025] Now, an exemplary illustration of the inventive scheduling concept will be presented, see Fig. 1. In this illustration, a computer system includes five processors: p1-p5 and three process groups g1, g2 and g3. The processors p1 and p2 constitute processor set s1. Processors p3, p4, and p5 are combined into processor set s2. Group g1 is assigned 3 shares and it is executing on the processor set s1. Process groups g2 and g3 are assigned 1 and 2 shares respectively, and they both execute on processor set s2.

[0026] The conventional fair share scheduler described above would allocate the system resources without reference to the processor sets. Specifically, it would assign $3/(3+1+2) = 3/6 = 50\%$ of the entire system resources to group g1. Group g2 would receive $1/(3+1+2) = 1/6 = 16.7\%$, while group g3 would get $2/(3+1+2) = 2/6 = 33.3\%$ of the system resources.

[0027] The inventive system, on the other hand, group g1 receives entire processor set s1, which is $40\% = 100\% * 2/5$ of the total system resources. Groups g2 and g3 jointly executing on processor set s2 receive $1/3 * 3/5 = 1/5 = 20\%$ and $2/3 * 3/5 = 2/5 = 40\%$ of the system resources, respectively. It should be noted that the percentage of the system resources assigned to process group g1 is independent of the number of the process groups and the load of the processor set s2.

[0028] According to another embodiment, a scheduler for fair share scheduling may be provided for allocating a percentage of system resources among process groups in a computer system. The computer system may include at least one central processing unit, said at least one central processing unit combined into at least one processor set. The central processing units may be located on one or a plurality of data processing devices, and may communicate with one another through internal communication links or external communication links.

[0029] The scheduler may be provided as an integral part of the computer system, e.g. in the form of a program or group of programs executed on at least one of the central processing units of the computer system or may be provided with at least one external device.

[0030] The scheduler for may comprise means for realizing the functionality outlined with respect to the previous embodiments. In particular, the scheduler may include means for assigning each of said process groups a number of shares for each or said at least one processor set. and means for allocating said system resources of each of said at least one processor set to each of

said process groups according to the number of shares assigned to said each of said process groups.

[0031] Further, the scheduler may include means for allocating said system resources of each of said at least one processor set based on a number of shares of all active groups within said each of said at least one processor set.

[0032] Moreover, the scheduler may include means for calculating said percentage of said system resources based on a ratio of the number of shares assigned to said each of said process groups to the a number of shares of all active groups within said each of said at least one processor set.

[0033] Each of said process groups may include only one process.

[0034] While the invention has been described herein with reference to preferred embodiments thereof, it will be readily apparent to persons of skill in the art that various modifications in form of detail can be made with respect thereto without departing from the spirit and scope of the invention as defined in and by the appended claims. For example, the present invention is not limited allocating CPU shares only. Generally, processes, users, or groups of users can be allocated shares relating to various system resources, such as CPU, disk or memory usage. In addition, shares need not be assigned to groups of users. According to the inventive concept, shares can be assigned to processes, users, or groups of users, or in any other manner. It will also be appreciated by those of skill in the art, that a number of processors in a processor set may not be an integer number.

[0035] In other words, the concept of processor set should be viewed, in the context of the present invention, as allocating shares of the total processor resources of a computer system among virtual computers (processor sets), that do not share processes, rather than simply combining resources of physical CPUs. Finally, each process group could be given a separate number of shares for each processor set.

[0036] Those of skill in the art will undoubtedly appreciate that the invention can be implemented on a wide variety of computer systems including, but not limited to, general purpose computers and special purpose computers such as network appliances. As well known in the art, a computer consists at least of a central processing unit, a memory unit, and an input/output interface. The aforementioned computer components can be arranged separately, or they can be combined together into a single unit. The computer memory unit may include a random access memory (RAM) and/or read only memory (ROM). The present invention can be implemented as a computer program embodied in any tangible storage medium, or loaded into the computer memory by any known means. As an alternative to implementing the present invention as a computer program, the present invention can be also embodied into an electronic circuit. This embodiment may provide an

improved performance characteristics.

[0037] According to another embodiment, a computer system may be provided comprising at least a central processing unit and a memory, said memory storing a program for allocating a percentage of system resources among process groups in a computer system, said computer system comprising at least one central processing unit, said at least one central processing unit combined into at least one processor set, said program comprising instructions for assigning each of said process groups a number of shares for each or said at least one processor set allocating said system resources of each of said at least one processor set to each of said process groups according to the number of shares assigned to said each of said process groups.

[0038] Further, it is noted that a computer-readable medium may be provided having a program embodied thereon, where the program is to make a computer or a system of data processing devices to execute functions or operations of the features and elements of the above described examples. A computer-readable medium can be a magnetic or optical or other tangible medium on which a program is recorded, but can also be a signal, e.g. analog or digital, electronic, magnetic or optical, in which the program is embodied for transmission. Further, a computer program product may be provided comprising the computer-readable medium.

Claims

1. A method for allocating a percentage of system resources among process groups in a computer system, said computer system comprising at least one central processing unit, said at least one central processing unit combined into at least one processor set, said method comprising:
 - a. assigning each of said process groups a number of shares for each or said at least one processor set;
 - b. allocating said system resources of each of said at least one processor set to each of said process groups according to the number of shares assigned to said each of said process groups.
2. The method of claim 1, wherein said system resources of each of said at least one processor set are allocated based on a number of shares of all active groups within said each of said at least one processor set.
3. The method of at least one of the claims 1 and 2, wherein said percentage of said system resources is calculated based on a ratio of the number of shares assigned to said each of said process

groups to the a number of shares of all active groups within said each of said at least one processor set.

4. The method of at least one of the claims 1 to 3, wherein each of said process groups includes only one process.
5. A computer readable medium embodying a program for allocating a percentage of system resources among process groups in a computer system, said computer system comprising at least one central processing unit, said at least one central processing unit combined into at least one processor set, said program comprising:
 - a. assigning each of said process groups a number of shares for each or said at least one processor set;
 - b. allocating said system resources of each of said at least one processor set to each of said process groups according to the number of shares assigned to said each of said process groups.
6. The computer readable medium of claim 5, wherein said system resources of each of said at least one processor set are allocated based on a number of shares of all active groups within said each of said at least one processor set.
7. The computer readable medium of at least one of the claims 5 and 6, wherein said percentage of said system resources is calculated based on a ratio of the number of shares assigned to said each of said process groups to the a number of shares of all active groups within said each of said at least one processor set.
8. The computer readable medium of at least one of the claims 5 to 8, wherein each of said process groups includes only one process.
9. A program having instructions adapted to make a computer carry out the method of at least one of the claims 1 - 4.
10. A scheduler for allocating a percentage of system resources among process groups in a computer system having at least one central processing unit, said at least one central processing unit combined into at least one processor set, the scheduler comprising:
 - means for assigning each of said process groups a number of shares for each or said at least one processor set; and

means for allocating said system resources of each of said at least one processor set to each of said process groups according to the number of shares assigned to said each of said process groups.

5

11. The scheduler of claim 10, including means for allocating said system resources of each of said at least one processor set based on a number of shares of all active groups within said each of said at least one processor set.

10

12. The scheduler of at least one of the claims 10 and 11, including means for calculating said percentage of said system resources based on a ratio of the number of shares assigned to said each of said process groups to the a number of shares of all active groups within said each of said at least one processor set.

15

20

13. The scheduler of at least one of the claims 10 to 12, wherein each of said process groups includes only one process.

25

30

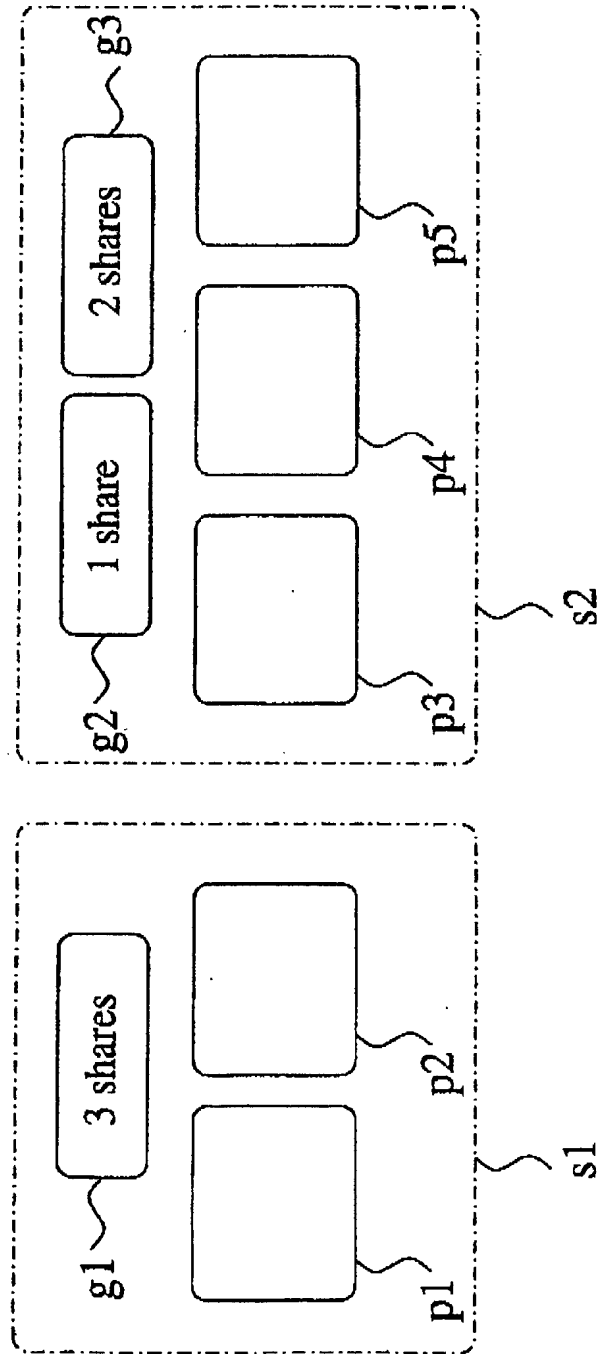
35

40

45

50

55



$P_{1-5} = 20\%$ (1/5th) of the system resources.

Fig. 1